

CLAIMS

What is claimed is:

1. A spam detection system comprising:
a component that identifies features relating to at least a portion of origination information of a message; and
a component that combines the features into useful pairs for use in connection with training a machine learning filter to facilitate detecting spam.
2. The system of claim 1, wherein each pair comprises at least one of the following:
at least one of a domain name and a host name in a MAIL FROM command;
at least one of a domain name and a host name in a HELO COMMAND;
at least one of an IP address and a subnet in a Received from header;
at least one of a domain name and a host name in a Display name;
at least one of a domain name and a host name in a Message From line;
and
at least one time zone in a last Received from header.
3. The system of claim 2, wherein the domain name is derived from the host name.
4. The system of claim 2, wherein the subnet comprises one or more IP addresses that share a first number of bits in common.
5. The system of claim 1, wherein a useful pair is any one of a domain name and a host name from a Message From and from a HELO command.
6. The system of claim 1, wherein a useful pair is a Display name domain name and host name and a Message From domain name and host name.

7. The system of claim 1, wherein a useful pair is any one of a domain name and a host name in a Message From and any one of a Received from IP address and subnet.

8. The system of claim 1, wherein a useful pair is a sender's alleged time zone and a Message From domain name.

9. The system of claim 1, wherein a useful pair comprises a sender's type of mailing software and any one of a domain name, host name and user name derived from one of an SMTP command and a message header.

10. The system of claim 1, wherein origination information comprises SMTP commands, the SMTP commands comprise a HELO command, a MAIL FROM command, and a DATA command.

11. The system of claim 10, wherein the DATA command comprises a Message From line, sender's alleged time zone, and sender's mailing software.

12. The system of claim 1, further comprising a component that applies one or more heuristics consistently to mail messages to obtain consistent feature pairing.

13. A spam detection system comprising:
a component that analyzes a portion of a message via searching for particular character sequences that are indicative of spam, wherein the particular sequences are not restricted to whole words; and
a component that generates features relating to the character sequences of any length.

14. The system of claim 13, wherein the component generates features for each run of characters up to a maximum character run length.

15. The system of claim 13, wherein the component generates features for substantially all character sequences up to some length n .

16. The system of claim 13, wherein the character sequences comprise at least one of letters, numbers, punctuation, symbols, and characters of foreign languages.

17. The system of claim 13, wherein the particular character sequences comprise at least one of random letters, symbols, and punctuation as chaff at any one of a beginning and end of at least one of a subject line of a message and a message body.

18. The system of claim 17, wherein random character sequences comprise character n -grams which are indicative of spam-like messages.

19. The system of claim 18, wherein the character n -grams are located in at least one of From address, subject line, text body, html body, and attachments.

20. The system of claim 18, wherein the character n -grams are position dependent.

21. The system of claim 13 for use with the messages comprising at least one of foreign language text, Unicode character types, and other character types not common to English

22. The system of claim 21, wherein the foreign language text comprises substantially non-space separated words.

23. The system of claim 22, wherein n -grams are used only for characters not typically separated by spaces.

24. The system of claim 13, further comprising a component that extracts character sequences obfuscated by punctuation using a pattern-match technique.
25. A spam detection system comprising:
a component that analyzes a portion of a message via searching for instances of a string of random characters that are indicative of the message being spam.
26. The system of claim 25, further comprising a component that generates features corresponding to the instances of random character strings to facilitate determining an entropy measurement for each string.
27. The system of claim 25, wherein the system measures a value correlated with entropy.
28. The system of claim 27, wherein a high value correlated with entropy is indicative of spam.
29. The system of claim 28, wherein the value correlated with entropy is the actual entropy $-\log_2 P(\text{abc} \dots \text{z})$
30. The system of claim 27, wherein the average entropy of a character string is used.
31. The system of claim 25, wherein the string of random characters is chaff.
32. The system of claim 27, wherein the relative entropy compares the entropy measurement at any one of a beginning and end of at least one of a subject line and message body with the entropy measurement at a middle of at least one of the subject line and message body.
33. A spam detection system comprising:

a component that analyzes substantially all features of a message header in connection with training a machine learning spam filter.

34. The system of claim 33, wherein the features of the message header comprise at least one of a presence and absence of at least one message header type, the message header types comprising X-Priority, mail software, and headers line for unsubscribing.

35. The system of claim 34, wherein the features of the message header further comprise content associated with at least one message header type.

36. The system of claim 33, further comprising:
a component that analyzes at least a portion of a message for images and related image information; and
a component that generates features relating to any one of the images and related image information.

37. The system of claim 36, wherein the image information comprises image size, image quantity, location of image, image dimensions, and image type.

38. The system of claim 36, wherein the image information comprises the presence of a first URL and a second URL such that the image is inside of a hyperlink.

39. The system of claim 38, wherein the message comprises a tag pattern having the form of </A.

40. The system of claim 36, wherein the features are used in connection with training a machine learning filter.

41. The system of claim 33, further comprising a component that analyzes a message for HTML attributes and location of HTML attributes as they appear in a tag pattern.

42. A method that facilitates generating features for use in spam detection comprising:
receiving at least one message;
parsing at least a portion of a message to generate one or more features;
combining at least two features into pairs, whereby each pair of features creates at least one additional feature, the features of each pair coinciding with one another; and
using the pairs of features to train a machine learning spam filter.

43. The method of claim 42, wherein the at least a portion of the message being parsed corresponds to origination information of the message.

44. The method of claim 42, wherein each pair comprises at least one of the following:
at least one of a domain name and a host name in a MAIL FROM command;
at least one of a domain name and a host name in a HELO-COMMAND;
at least one of an IP address and a subnet in a Received from header;
at least one of a domain name and a host name in a Display name;
at least one of a domain name and a host name in a Message From line;
and
at least one time zone in a last Received from header.

45. The method of claim 44, wherein the domain name is derived from the host name.

46. The method of claim 42, wherein the pair of features is a Display name domain name and host name and a Message From domain name and host name.

47. The method of claim 42, wherein a useful pair is any one of a domain name and a host name from a Message From and from a HELO command.

48. The method of claim 42, wherein the pair of features is any one of a domain name and a host name in a Message From and any one of a Received from IP address and subnet.

49. The method of claim 42, wherein the pair of features is a sender's alleged time zone and a Message From domain name.

50. The method of claim 42, wherein the pair of features comprises a sender's type of mailing software and any one of a domain name, host name and display name derived from one of an SMTP command and a message header.

51. The method of claim 42, further comprising selecting one or more most useful pairs of features to train the machine learning filter.

52. The method of claim 42, further comprising employing the machine learning filter after it is trained to detect spam by performing the following:

receiving new messages;

generating pairs of features based on origination information in the messages;

passing the pairs of features through the machine learning filter; and

obtaining a verdict as to whether at least one pair of features indicates that the message is more likely to be spam.

53. A method that facilitates generating features for use in spam detection comprising:

receiving one or more messages;
walking through at least a portion of the message to create features for each run of characters of any run length; and
training a machine learning filter using at least a portion of the created features.

54. The method of claim 53, further comprising generating features relating to a position of at least one run of characters.

55. The method of claim 54, wherein the position comprises any one of a beginning of a message body, an end of a message body, a middle of a message body, a beginning of a subject line, an end of a subject line, and a middle of a subject line.

56. The method of claim 53, wherein the features are created for a run of characters up to length n .

57. The method of claim 53, wherein the features are created for sub-lengths of runs of characters.

58. The method of claim 53, wherein the run of characters comprise character n -grams.

59. The method of claim 53, further comprising calculating an entropy of one or more run of characters and employing the calculated entropy as a feature in connection with training a spam filter.

60. The method of claim 59, wherein the entropy is at least one of high entropy, average entropy, and relative entropy.

61. The method of claim 60, wherein the average entropy is the entropy per character of a particular run of characters.

62. The method of claim 60, wherein the relative entropy is a comparison of the entropy of a particular run of characters at a first location relative to the entropy of a particular run of characters at a second location of the message.

63. The method of claim 62, wherein the first and second locations comprise a beginning of a subject line, a middle of a subject line, and an end of a subject line, whereby the first location is not the same as the second location when determining the relative entropy for any given run of characters.

64. The method of claim 62, wherein the first and second locations comprise a beginning of a message, a middle of a message, and an end of a message, whereby the first location is not the same as the second location when determining the relative entropy for any given run of characters.

65. The method of claim 53, further comprising employing the machine learning filter after it is trained to detect spam by performing the following:

- receiving new messages;
- generating features based at least one of runs of characters and entropy determinations of runs of characters in the messages;
- passing the features through the machine learning filter; and
- obtaining a verdict as to whether the features indicate that the message is more likely to be spam.

66. A method that facilitates generating features for use in spam detection comprising:

- receiving one or more messages;
- analyzing substantially all features of a message header; and
- training a machine learning filter using the analyzed features.

67. The method of claim 66, further comprising analyzing substantially all features based on image information in the message.

68. A computer readable medium comprising the method of claim 42.

69. A computer readable medium comprising the method of claim 53.

70. A computer-readable medium having stored thereon the following computer executable components:

a component that identifies features relating to at least a portion of origination information of a message; and

a component that combines the features into useful pairs for use in connection with training a machine learning filter to facilitate detecting spam.

71. The computer readable medium of claim 70, further comprising:

a component that analyzes a portion of a message via searching for particular character sequences that are indicative of spam, wherein the particular sequences are not restricted to whole words; and

a component that generates features relating to the character sequences of any length.

72. The computer readable medium of claim 70, further comprising:

a component that analyzes a portion of a message via searching for instances of a string of random characters that are indicative of the message being spam.

73. A system that facilitates generating features for use in spam detection comprising:

a means for receiving at least one message;

a means for parsing at least a portion of a message to generate one or more features;

a means for combining at least two features into pairs, whereby each pair of features creates at least one additional feature, the features of each pair coinciding with one another; and

a means for using the pairs of features to train a machine learning spam filter.

74. A system that facilitates generating features for use in spam detection comprising:

a means for receiving one or more messages;

a means for walking through at least a portion of the message to create features for each run of characters of any run length; and

a means for training a machine learning filter using at least a portion of the created features.

75. The system of claim 74, further comprising calculating an entropy of one or more run of characters and employing the calculated entropy as a feature in connection with training a spam filter.